

Internal Assessment Test 1 – Nov 2024

|                                  |  |           |         |            |           |            |           |      |     |     |
|----------------------------------|--|-----------|---------|------------|-----------|------------|-----------|------|-----|-----|
| Sub:                             | INFORMATION RETRIEVAL  |           |         |            | Sub Code: | BAI515B    | Branch:   | AIML |     |     |
| Date:                            | 06 / 11/2024   | Duration: | 90 mins | Max Marks: | 50        | Sem / Sec: | V / A,B,C |      |     | OBE |
| <u>Answer Any of 5 Questions</u> |  |           |         |            |           |            | MARKS     | CO   | RBT |     |
| 1                                | Explain the Process of Information Retrieval and the components involved in it with a neat architecture.   |           |         |            |           | [10]       | CO1       | L2   |     |     |
| 2 (a)                            | Define the Vector Model and the advantages of the Vector Model?  |           |         |            |           | [05]       | CO2       | L2   |     |     |
| (b)                              | Can the TF-IDF weight of a term in a document exceed? why?   |           |         |            |           | [05]       | CO2       | L2   |     |     |
| 3                                | Consider the following Documents<br>Query- Obama health Plan<br>D1: Obama rejects allegations about his own bad health<br>D2: The Plan is to visit Obama<br>D3:Obama raises concerns with US health plan reforms<br>Estimate the probability that above Documents are Relevant to the Query. |           |         |            |           | [10]       | CO2       | L3   |     |     |
| 4 (a)                            | Explain Receiver Operating Characteristics and Benefits of ROC   |           |         |            |           | [05]       | CO3       | L2   |     |     |
| (b)                              | Consider the Two texts,"Tom and Jerry are friends" and "Jack and Tom are friends". Calculate the Cosine similarity for these two texts?  |           |         |            |           | [05]       | CO2       | L3   |     |     |
| 5.a)                             | Explain the Types of Text Compression Techniques.  |           |         |            |           | [05]       | CO3       | L2   |     |     |
| b)                               | How Does the Large amount of Information available in Web affect information retrieval system Implementation?  |           |         |            |           | [05]       | CO1       | L2   |     |     |
| 6                                | If an IR System returns 6 relevant Documents and 10 non relevant documents.there are Total of 20 relevant Documents in the collection.calculate the Precision and Recall of the system on this search?   |           |         |            |           | [10]       | CO3       | L3   |     |     |

CI CCI HOD-AIML

---

---

## Internal Assessment Test 1 – Nov 2024

| Sub:                             | INFORMATION RETRIEVAL  | Sub Code:  | BAI515B   | Branch:    | AIML  |    |     |
|----------------------------------|--|------------|-----------|------------|-------|----|-----|
| Date:                            | 06 / 11/2024   | Duration:  | 90 mins   | Max Marks: | 50    |    |     |
|                                  |  | Sem / Sec: | V / A,B,C |            | OBE   |    |     |
| <u>Answer Any of 5 Questions</u> |  |            |           |            | MARKS | CO | RBT |
| 1                                | Explain the Process of Information Retrieval and the components involved in it with a neat architecture.<br>Definition-4<br>Drawing- 4<br>Explanation-2  | [10]       | CO1       | L2         |       |    |     |
| 2 (a)                            | Define the Vector Model and the advantages of the Vector Model?<br>Definition & Examples-3<br>advantages-2   | [05]       | CO2       | L2         |       |    |     |
| (b)                              | Can the TF-IDF weight of a term in a document exceed? why?<br>Definition & Examples-3<br>why TF-IDF-2  | [05]       | CO2       | L2         |       |    |     |
| 3                                | Consider the following Documents<br>Query- Obama health Plan<br>D1: Obama rejects allegations about his own bad health<br>D2: The Plan is to visit Obama<br>D3:Obama raises concerns with US health plan reforms<br>Estimate the probability that above Documents are Relevant to the Query.<br><br>Stop wor & Stemming-2 M<br>Query Comparison with Documents- 4M<br>Ranking-4M | [10]       | CO2       | L3         |       |    |     |
| 4 (a)                            | Explain Receiver Operating Characteristics and Benefits of ROC<br>Definition-3M<br>Benefits-2M   | [05]       | CO3       | L2         |       |    |     |
| (b)                              | Consider the Two texts,"Tom and Jerry are friends" and "Jack and Tom are friends". Calculate the Cosine similarity for these two texts?<br>Definition-3M<br>Steps-2M   | [05]       | CO2       | L3         |       |    |     |
| 5.a)                             | Explain the Types of Text Compression Techniques.<br>Definition-3M<br>Types-2M   | [05]       | CO3       | L2         |       |    |     |
| b)                               | How Does the Large amount of Information available in Web affect information retrieval system Implementation?<br><br>Definition-3M<br>Types of Web-2M  | [05]       | CO1       | L2         |       |    |     |
| 6                                | If an IR System returns 6 relevant Documents and 10 non relevant documents.there are Total of 20 relevant Documents in the collection.calculate the Precision and Recall of the system on this search?<br><br>Definition-5M<br>Steps-2M<br>Precision & recall-3M   | [10]       | CO3       | L3         |       |    |     |





USN 

|  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|
|  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|



Internal Assessment Test 1 – Nov 2024

|       |                       |           |         |            |    |            |             |         |      |     |     |
|-------|-----------------------|-----------|---------|------------|----|------------|-------------|---------|------|-----|-----|
| Sub:  | INFORMATION RETRIEVAL |           |         |            |    | Sub Code:  | BAI515B     | Branch: | AIML |     |     |
| Date: | 06 / 11/2024          | Duration: | 90 mins | Max Marks: | 50 | Sem / Sec: | V / A, B, C |         |      | OBE |     |
|       |                       |           |         |            |    |            |             |         | M    | CO  | RBT |
|       |                       |           |         |            |    |            |             |         | A    |     |     |
|       |                       |           |         |            |    |            |             |         | R    |     |     |
|       |                       |           |         |            |    |            |             |         | K    |     |     |
|       |                       |           |         |            |    |            |             |         | S    |     |     |

Answer Any of 5 Questions

1. Explain the Process of Information Retrieval and the components involved in it with a neat architecture. [10] CO1 L2

A. **Information Retrieval (IR)** can be defined as a software program that deals with the organization, storage, retrieval, and evaluation of information from document repositories, particularly textual information. Information Retrieval is the activity of obtaining material that can usually be documented on an unstructured nature i.e. usually text which satisfies an information need from within large collections which is stored on computers. For example, Information Retrieval can be when a user enters a query into the system.

**Components of Information Retrieval (IR)**

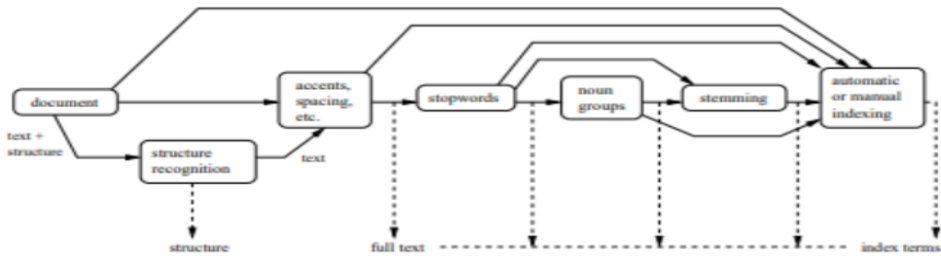


Figure 1.2 Logical view of a document: from full text to a set of index terms.

**Document Processing:** This is the first step where raw data, like text documents, are prepared for indexing. The goal here is to transform unstructured data into a structured format that is easier to work with.

Document processing includes several sub-steps:

- **Tokenization:** This involves breaking down text into smaller parts, called "tokens." Each token typically represents a word. For example, the sentence "The cat is cute" is split into tokens like "The," "cat," "is," and "cute." Tokenization makes it easier to work with words as individual pieces of data.
- **Stemming/Lemmatization:** Words are often reduced to their base or root form. For instance, words like "running," "ran," and "runs" can all be reduced to "run." This helps group similar words and improves the matching between documents and queries.
- **Stop Words Removal:** Common words such as "the," "is," and "and" don't add much meaning to searches, so they're removed. This process helps focus on the essential words that are more likely to determine the document's content.

**Indexing:** Once the documents are processed, an index is created. Think of indexing like creating a library catalogue where each word points to the documents containing it. An inverted index is commonly used, which maps each word (term) to the documents (IDs) where it appears. This makes searching fast, as we can quickly look up terms and find relevant documents without scanning everything.

**Query Processing:** When a user submits a query, it's processed to match the indexing structure. This makes the query ready for comparison against the indexed data.

Query processing often includes:

- **Query Tokenization and Expansion:** Similar to document tokenization, the query is split into tokens, and sometimes expanded with synonyms or related words. For example, a search for "car" might include "automobile" or "vehicle" to cover related terms and improve search results.
- **Relevance Feedback:** After an initial set of results, relevance feedback can refine the search. For example, if the user clicks on certain documents, the system may adjust future results to include more similar documents.

**Ranking and Retrieval:** Algorithms rank documents based on how closely they match the query. Scoring methods like TF-IDF (Term Frequency-Inverse Document Frequency) or BM25 are commonly used. These methods assign scores to documents based on term frequency, rarity, and other factors. Higher scores mean higher relevance, so the system retrieves the most relevant documents first.

**Evaluation:** This final step assesses the IR system's performance using metrics such as precision (how many retrieved documents are relevant), recall (how many relevant documents were retrieved out of all relevant ones), and the F1-score (a balance between precision and recall). Evaluation helps ensure that the system provides accurate and useful results.

|       |  |      |     |    |
|-------|--|------|-----|----|
| 2 (a) | <p>Define the Vector Model and the advantages of Vector Model?</p> <p>In the <a href="#">Vector Space</a> Model (VSM), each document or query is a N-dimensional vector where N is the number of distinct terms over all the documents and queries. The i-th index of a vector contains the score of the i-th term for that vector.</p> <p>The main score functions are based on: Term-Frequency (tf) and Inverse-Document-Frequency(idf).</p> <p><b>Term Frequency (TF)</b></p> <p>The <b>Term Frequency</b> <math>tf_{i,j}</math> measures the frequency of the i-th term in the j-th document. It is calculated by dividing the number of occurrences of term <math>t_i</math> in document <math>j</math> by the total number of terms in document <math>j</math>:</p> $tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$ <ul style="list-style-type: none"> <li>• <math>n_{i,j}</math> is the number of occurrences of term <math>i</math> in document <math>j</math>,</li> <li>• <math>\sum_k n_{k,j}</math> is the total number of occurrences of all terms in document <math>j</math>.</li> </ul> <p><b>Inverse Document Frequency (IDF)</b></p> <p>The <b>Inverse Document Frequency</b> <math>idf_i</math> evaluates the importance of term <math>i</math> across all documents. Rare terms are given higher weights as they are considered more specific to a document. It is computed as:</p> $idf_i = \log \frac{ D }{ \{d : t_i \in d\} }$ <p>where:</p> <ul style="list-style-type: none"> <li>• <math> D </math> is the total number of documents,</li> <li>• <math> \{d : t_i \in d\} </math> is the number of documents containing term <math>i</math>.</li> </ul> <p><b>Cosine Similarity</b></p> <p>To compute the similarity between two vectors <math>a</math> and <math>b</math> (representing document-query or document-document pairs), we use <b>Cosine Similarity</b>. The cosine of the angle between vectors <math>a</math> and <math>b</math> is calculated as:</p> $\cos(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{\ \mathbf{a}\  \ \mathbf{b}\ } = \frac{\sum_{i=1}^n \mathbf{a}_i \mathbf{b}_i}{\sqrt{\sum_{i=1}^n (\mathbf{a}_i)^2} \sqrt{\sum_{i=1}^n (\mathbf{b}_i)^2}}$ <p>where:</p> <ul style="list-style-type: none"> <li>• <math>\mathbf{a} \cdot \mathbf{b}</math> is the dot product of vectors <math>\mathbf{a}</math> and <math>\mathbf{b}</math>,</li> <li>• <math>\ \mathbf{a}\ </math> and <math>\ \mathbf{b}\ </math> are the magnitudes (norms) of vectors <math>\mathbf{a}</math> and <math>\mathbf{b}</math>.</li> </ul> <p><b>Advantages:</b></p> <ul style="list-style-type: none"> <li>• The retrieval performance is improved by its term-weighting method.</li> <li>• With the help of its partial matching approach, documents that roughly match the query conditions can be retrieved.</li> <li>• The documents are sorted using its cosine ranking formula based on how similar they are to the query.</li> </ul> | [05] | CO2 | L2 |
|-------|--|------|-----|----|

|       |   |      |     |    |
|-------|---|------|-----|----|
| 2 (b) | <p>Can the TF-IDF weight of a term in a document exceed? why?</p> <p><b>1. Term Frequency (TF):</b></p> <ul style="list-style-type: none"> <li>• <b>Definition:</b> Term Frequency measures how frequently a term appears in a document. The more frequently a term appears, the more important it is assumed to be for that document.</li> <li>• <b>Formula:</b></li> </ul> $TF = \frac{\text{Number of times a term appears in a document}}{\text{Total number of terms in the document}}$ <ul style="list-style-type: none"> <li>• <b>Purpose:</b> TF helps determine the relevance of a term within a single document. A higher term frequency indicates that the term is more significant in that specific document.</li> </ul> <p><b>2. Inverse Document Frequency (IDF):</b></p> <ul style="list-style-type: none"> <li>• <b>Definition:</b> Inverse Document Frequency measures how unique or rare a term is across all documents in a collection. It helps reduce the weight of common terms (like "the" or "is") that appear in many documents and increases the weight of rare terms.</li> <li>• <b>Formula:</b></li> </ul> $IDF = \log \left( \frac{\text{Total number of documents}}{\text{Number of documents containing the term}} \right)$ <ul style="list-style-type: none"> <li>• <b>Purpose:</b> IDF gives more weight to terms that are specific to fewer documents, which often makes them more relevant to a specific topic.</li> </ul> <p>Yes, the <b>TF-IDF</b> weight of a term in a document can, in theory, be <b>greater than 1</b>.</p> <p><b>1. High Term Frequency (TF) Contribution:</b></p> <ul style="list-style-type: none"> <li>• If a term appears frequently within a document (high term frequency), the TF component (e.g., <math>tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}</math>) can increase significantly, contributing to a larger overall TF-IDF value.</li> </ul> <p><b>2. Logarithmic IDF Function:</b></p> <ul style="list-style-type: none"> <li>• The IDF component is typically a logarithmic function, such as <math>idf_i = \log \frac{ D }{ \{d: t_i \in d\} }</math>. If a term is rare across the document corpus, the IDF value can be large, especially if the document collection is vast, resulting in a higher TF-IDF score.</li> </ul> <p><b>3. No Normalization by Default:</b></p> <ul style="list-style-type: none"> <li>• In standard TF-IDF calculation, there's no explicit normalization to constrain weights between 0 and 1, so TF-IDF values can exceed 1 based on the term's relative frequency and rarity.</li> </ul> | [05] | CO2 | L2 |
|-------|---|------|-----|----|



3

Consider the following Documents

Query- Obama health Plan

D1: Obama rejects allegations about his own bad health

D2: The Plan is to visit Obama

D3: Obama raises concerns with US health plan reforms

Estimate the probability that above Documents are Relevant to the Query.

The keywords from the Query are: “Obama”, “health”, “Plan”.

Here is the **Contingency table** the given documents and query:

| Document     | Obama    | Health     | Plan       | Total    |
|--------------|----------|------------|------------|----------|
| Doc 1        | 3/3 = 1  | 2/3 = 0.67 | 0/3 = 0    | 2        |
| Doc 2        | 3/3 = 1  | 0/3 = 0    | 2/3 = 0.67 | 2        |
| Doc 3        | 3/3 = 1  | 2/3 = 0.67 | 2/3 = 0.67 | 3        |
| <b>Total</b> | <b>3</b> | <b>2</b>   | <b>2</b>   | <b>7</b> |

The probability for the keywords of Query to exist in the documents is given below

$$\text{Probability of Doc 1} = \frac{3 \times 1 + 2 \times 0.67 + 0 \times 2}{7} = 0.62$$

$$\text{Probability of Doc 2} = \frac{3 \times 1 + 0 \times 0.67 + 2 \times 0.67}{7} = 0.62$$

$$\text{Probability of Doc 3} = \frac{3 \times 1 + 2 \times 0.67 + 2 \times 0.67}{7} = 0.9152$$

From the above data it is clear that **Doc3** is having more probability of existence with a probability of **0.9152 (91.52 %)**. So, the order of relevance is **Doc3 > {Doc1, Doc2}** [Since Doc2 and Doc1 are having equal probabilities they are placed in the same set]

[10]

CO2

L3

4 (a) Explain Receiver Operating Characteristics and Benefits of ROC

[05]

CO3

L2

**Receiver Operating Characteristics (ROC)** is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold varies. It's widely used to evaluate models, particularly in binary classification problems.

1. **True Positive Rate (TPR)** (also known as Sensitivity or Recall):

- The proportion of actual positives correctly identified by the model.
- Calculated as  $TPR = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

2. **False Positive Rate (FPR):**

- The proportion of actual negatives incorrectly identified as positives.
- Calculated as  $FPR = \frac{\text{False Positives}}{\text{False Positives} + \text{True Negatives}}$

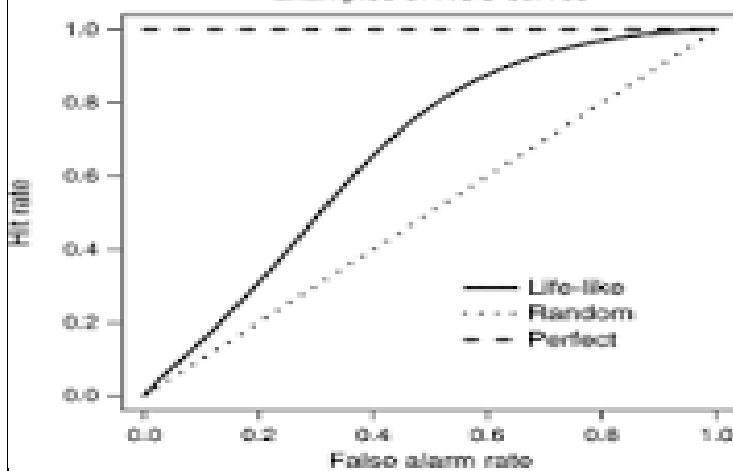
3. **ROC Curve:**

- The ROC curve is a plot of the TPR against the FPR at various threshold settings. Each point on the curve represents a different threshold for classifying a data point as positive or negative.
- A model with good performance will have a curve that moves toward the top-left corner, showing a high TPR and a low FPR.

4. **Area Under the Curve (AUC):**

- The AUC-ROC score is the area under the ROC curve. An AUC close to 1 indicates a strong classifier, while an AUC of 0.5 indicates no better than random guessing.

Examples of ROC curves



ROC curves are used in many fields, including:

- **Medical diagnostics:** ROC curves are a well-known tool for evaluating the accuracy of diagnostic tests.
- **Epidemiology:** ROC curves can be used in epidemiology.
- **Radiology:** ROC curves can be used in radiology.
- **Bioinformatics:** ROC curves can be used in bioinformatics.
- **Stock market:** ROC curves can be used in the stock market.
- **Fruit tree survival:** ROC curves can be used to predict fruit tree survival.
- **Sports:** ROC curves can be used in sports.

4 (b) Consider the Two texts, "Tom and Jerry are friends" and "Jack and Tom are friends". Calculate the Cosine similarity for these two texts? [05] CO2 L3

Step 1: Tokenize the text<sub>1</sub>

Text<sub>1</sub>: "Tom and Jerry are friends"

Tokens: ["Tom", "and", "Jerry", "are", "friends"]

Step 2: tokenize the text<sub>2</sub>

Text<sub>2</sub>: "Jack and Tom are friends"

Tokens: ["Jack", "and", "Tom", "are", "friends"]

By combining Tokens

Token: ["Tom", "and", "Jerry", "are", "Friends", "Jack"]

Step 3: Prepare term frequency vector

for text<sub>1</sub>:  
 $v_1 = [1, 1, 1, 1, 1, 0]$

for text<sub>2</sub>:  
 $v_2 = [1, 1, 0, 1, 1, 1]$

Step 4:

$$\text{Cosine Similarity} = \frac{A \cdot B}{|A||B|}$$

$$= \frac{(1 \times 1) + (1 \times 1) + (1 \times 0) + (1 \times 1) + (1 \times 1) + (0 \times 1)}{\left(\sqrt{1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 0}\right) \left(\sqrt{1^2 + 1^2 + 0^2 + 1^2 + 1^2 + 1^2}\right)}$$

$$= \frac{4}{5} \approx 0.8$$

$$\boxed{\text{Cosine Similarity} = 0.8}$$

|      |  |      |     |    |
|------|--|------|-----|----|
| 5.a) | <p>Explain the Types of Text Compression Techniques.</p> <p>In <b>Information Retrieval (IR)</b>, text compression techniques are essential for reducing storage requirements and improving retrieval efficiency. Text compression can be divided into <b>two main categories: lossless</b> and <b>lossy</b> compression. Lossless compression maintains the exact original text, while lossy compression allows some information to be discarded to achieve higher compression. Below are the primary types of text compression techniques used in IR:</p> <hr/> <p><b>1. Statistical Compression</b><br/> Statistical methods rely on analyzing the frequency of characters or patterns in the text and encoding more frequent items with shorter codes.</p> <ul style="list-style-type: none"> <li>• <b>Huffman Coding:</b> <ul style="list-style-type: none"> <li>○ This is a widely-used lossless technique that builds a binary tree based on character frequency. More frequent characters are given shorter binary codes, and less frequent characters are given longer codes.</li> </ul> </li> <li>• <b>Arithmetic Coding:</b> <ul style="list-style-type: none"> <li>○ It encodes the entire message as a single number within a range by recursively subdividing intervals based on symbol probabilities. It is highly efficient and can offer better compression rates than Huffman coding.</li> </ul> </li> <li>• <b>Shannon-Fano Coding:</b> <ul style="list-style-type: none"> <li>○ Similar to Huffman coding, Shannon-Fano coding assigns shorter codes to more frequent symbols. It's not as optimal as Huffman coding but is simpler.</li> </ul> </li> </ul> <hr/> <p><b>2. Dictionary-Based Compression</b><br/> These methods replace common phrases or words with shorter codes based on a dictionary.</p> <ul style="list-style-type: none"> <li>• <b>Lempel-Ziv-Welch (LZW):</b> <ul style="list-style-type: none"> <li>○ LZW builds a dictionary of substrings dynamically as the text is read. It replaces repeated patterns with shorter codes, improving compression for texts with many repeated phrases.</li> </ul> </li> <li>• <b>Ziv-Lempel (LZ77 and LZ78):</b> <ul style="list-style-type: none"> <li>○ Both LZ77 and LZ78 use sliding windows to match strings with previously seen sequences. LZ77 references previous occurrences of a string within a defined window, while LZ78 builds a dictionary as it processes text.</li> </ul> </li> </ul> <hr/> <p><b>3. Transform-Based Compression</b><br/> These techniques transform the text to improve compression efficiency by rearranging it into a form that's easier to encode.</p> <ul style="list-style-type: none"> <li>• <b>Burrows-Wheeler Transform (BWT):</b> <ul style="list-style-type: none"> <li>○ BWT rearranges the text so that similar characters are grouped together, which enhances the efficiency of other compression techniques like Run-Length Encoding (RLE) or Huffman Coding. This is the foundation for compression algorithms like <b>bzip2</b>.</li> </ul> </li> </ul> <hr/> <p><b>4. Run-Length Encoding (RLE)</b></p> <ul style="list-style-type: none"> <li>• <b>Run-Length Encoding:</b> <ul style="list-style-type: none"> <li>○ RLE is a straightforward technique where sequences of repeated characters are stored as a single character followed by a count. For instance, "aaaabbbb" would be encoded as "a4b4". It works well for texts with many repeated characters or patterns.</li> </ul> </li> </ul> <hr/> <p><b>5. Hybrid Methods</b></p> <ul style="list-style-type: none"> <li>• Some modern compression algorithms combine multiple techniques to maximize compression. For example, the <b>DEFLATE</b> algorithm used in ZIP files combines LZ77 (dictionary-based) and Huffman coding (statistical) for enhanced efficiency.</li> </ul> <hr/> <p><b>6. Lossy Compression Techniques</b></p> <ul style="list-style-type: none"> <li>• While rarely used for general text due to the need for exact data recovery, lossy techniques can be applied to certain types of IR data, like summarization or topic modeling, where approximate data is acceptable. Techniques like <b>vector quantization</b> and <b>Latent Semantic Analysis (LSA)</b> are sometimes used in this context.</li> </ul> | [05] | CO3 | L2 |
|------|--|------|-----|----|

|      |  |      |     |    |
|------|--|------|-----|----|
| 5 b) | <p>How Does the Large amount of Information available in Web affect information retrieval system Implementation?</p> <p>A <b>Web Information Retrieval System</b> is designed to gather, organize, index, and retrieve relevant information from the vast, dynamic content on the internet. Unlike traditional information retrieval systems that work with closed document collections, web IR systems must manage vast data, handle various document formats, and adapt to constantly changing content. Web IR powers search engines like Google, Bing, and others, focusing on speed, relevance, and user experience.</p> <p><b>Key Components of a Web Information Retrieval System</b></p> <ol style="list-style-type: none"> <li>1. <b>Web Crawling:</b> <ul style="list-style-type: none"> <li>○ Web crawlers (or spiders) navigate the internet to discover and retrieve web pages. They follow links to build a comprehensive and up-to-date index of available content, periodically revisiting pages to keep the index current.</li> </ul> </li> <li>2. <b>Indexing:</b> <ul style="list-style-type: none"> <li>○ After gathering web pages, the system indexes the content, which involves tokenizing text, removing stop words, stemming, and creating an <b>inverted index</b>. This index allows for rapid searching by mapping each term to the documents in which it appears.</li> </ul> </li> <li>3. <b>Document Representation and Metadata Extraction:</b> <ul style="list-style-type: none"> <li>○ Each web page is represented using document vectors (e.g., with TF-IDF or BM25 weighting) that quantify term relevance. Metadata, such as page titles, URLs, and tags, is also extracted to enhance retrieval quality.</li> </ul> </li> <li>4. <b>Query Processing and Ranking:</b> <ul style="list-style-type: none"> <li>○ User queries are analysed, tokenized, and possibly expanded to match indexed documents effectively. The IR system ranks documents using relevance-based algorithms (e.g., cosine similarity, BM25), link-based ranking (e.g., PageRank), and often considers user intent, context, and personalization factors.</li> </ul> </li> <li>5. <b>Relevance Feedback and Personalization:</b> <ul style="list-style-type: none"> <li>○ Based on user interactions like clicks and time spent on a page, the system adjusts rankings and tailors' future queries to individual preferences, improving relevance and user satisfaction.</li> </ul> </li> <li>6. <b>User Interface and Experience:</b> <ul style="list-style-type: none"> <li>○ A user-friendly interface displays search results, query suggestions, filters, and other interactive features. This UI is critical for effective search experience, as it influences how users interact with results and perceive relevance.</li> </ul> </li> </ol> <p>The vast amount of information available on the web significantly impacts the <b>implementation</b> and <b>effectiveness</b> of information retrieval (IR) systems. Here are the key challenges and considerations:</p> <ol style="list-style-type: none"> <li>1. <b>Scalability and Storage Requirements</b> <ul style="list-style-type: none"> <li>● <b>Challenge:</b> Web data is constantly growing, so IR systems must handle and store vast amounts of information.</li> <li>● <b>Solution:</b> Efficient indexing, distributed storage systems, and scalable infrastructure (e.g., cloud storage) are required to manage this growth.</li> </ul> </li> <li>2. <b>Speed and Latency in Retrieval</b> <ul style="list-style-type: none"> <li>● <b>Challenge:</b> Retrieving relevant documents quickly from a large dataset can lead to high latency.</li> <li>● <b>Solution:</b> Optimized indexing (e.g., inverted indexes) and caching are essential to reduce response times and improve user experience.</li> </ul> </li> <li>3. <b>Relevance and Precision of Results</b> <ul style="list-style-type: none"> <li>● <b>Challenge:</b> Large-scale data leads to diverse content, making it challenging to retrieve only the most relevant results.</li> <li>● <b>Solution:</b> Advanced ranking algorithms, personalized search, and user behavior analysis can help IR systems surface more relevant results.</li> </ul> </li> <li>4. <b>Handling Data Diversity and Quality</b> <ul style="list-style-type: none"> <li>● <b>Challenge:</b> Web content varies widely in language, structure, quality, and relevance.</li> <li>● <b>Solution:</b> Preprocessing steps like language detection, stopword removal, and content filtering are necessary to handle such diversity.</li> </ul> </li> </ol> | [05] | CO1 | L2 |
|------|--|------|-----|----|

|            |  |      |     |    |
|------------|--|------|-----|----|
| 6<br><br>A | <p>If an IR System returns 6 relevant Documents and 10 non relevant documents.there are Total of 20 relevant Documents in the collection.calculate the Precision and Recall of the system on this search?</p> <p><b>Precision and Recall Calculation for Information Retrieval System</b><br/>In the context of evaluating the performance of an Information Retrieval (IR) system, <b>Precision</b> and <b>Recall</b> are two important metrics. Let's calculate these metrics based on the given data:</p> <hr/> <p><b>Given:</b></p> <ul style="list-style-type: none"> <li>• The IR system returns <b>6 relevant documents</b>.</li> <li>• The IR system returns <b>10 non-relevant documents</b>.</li> <li>• The total number of <b>relevant documents in the collection is 20</b>.</li> </ul> <p>The formulae to calculate Precision and Recall are:</p> <p><b>1. Precision</b><br/>Precision is the measure of how many of the documents retrieved by the system are actually relevant.</p> $\text{Precision} = \frac{\text{Number of Relevant Documents Retrieved}}{\text{Total Documents Retrieved}}$ <p>Where:</p> <ul style="list-style-type: none"> <li>• The number of relevant documents retrieved = 6.</li> <li>• The total number of documents retrieved = 16 (6 relevant + 10 non-relevant).</li> </ul> <p>Substituting the values:</p> $\text{Precision} = \frac{6}{16} = 0.375$ <p>Thus, Precision = 0.375 or 37.5%.</p> <p>This means that 37.5% of the documents retrieved by the IR system are actually relevant.</p> <hr/> <p><b>2. Recall</b><br/>Recall is the measure of how many of the total relevant documents in the collection were successfully retrieved by the IR system.</p> $\text{Recall} = \frac{\text{Number of Relevant Documents Retrieved}}{\text{Total Relevant Documents in Collection}}$ <p>Where:</p> <ul style="list-style-type: none"> <li>• The number of relevant documents retrieved = 6.</li> <li>• The total number of relevant documents in the collection = 20.</li> </ul> <p>Substituting the values:</p> $\text{Recall} = \frac{6}{20} = 0.3$ <p>Thus, Recall = 0.3 or 30%.</p> | [10] | CO3 | L3 |
|------------|--|------|-----|----|

CI

CCI

HOD-AIML